



Mixed-norm Regularization for Brain Decoding

Rémi Flamary, Nisrine Jrad, Ronald Phlypo, Marco Congedo, Alain Rakotomamonjy

► To cite this version:

Rémi Flamary, Nisrine Jrad, Ronald Phlypo, Marco Congedo, Alain Rakotomamonjy. Mixed-norm Regularization for Brain Decoding. Computational and Mathematical Methods in Medicine, 2014, 2014, pp.ID 317056. 10.1155/2014/317056 . hal-00708243v2

HAL Id: hal-00708243

<https://hal.science/hal-00708243v2>

Submitted on 14 Mar 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Mixed-norm Regularization for Brain Decoding

R. Flamary^{a,*}, N. Jrad^c, R. Phlypo^c, M. Congedo^c, A. Rakotomamonjy^{b,*}

^a*Laboratoire Lagrange, UMR7293, Université de Nice, 00006 Nice, France*

^b*LITIS, EA 4108 - INSA / Université de Rouen, 76000 Rouen, France*

^c*Gipsa Lab, Domaine Universitaire BP 46, 38402 Saint Martin d'Hères cedex, France*

Abstract

This work investigates the use of mixed-norm regularization for sensor selection in Event-Related Potential (ERP) based Brain-Computer Interfaces (BCI). The classification problem is cast as a discriminative optimization framework where sensor selection is induced through the use of mixed-norms. This framework is extended to the multi-task learning situation where several similar classification tasks related to different subjects are learned simultaneously. In this case, multi-task learning helps in leveraging data scarcity issue yielding to more robust classifiers. For this purpose, we have introduced a regularizer that induces both sensor selection and classifier similarities. The different regularization approaches are compared on three ERP datasets showing the interest of mixed-norm regularization in terms of sensor selection. The multi-task approaches are evaluated when a small number of learning examples are available yielding to significant performance improvements especially for subjects performing poorly.

Keywords: Brain Computer Interface, Support Vector Machines, Sensor selection, EEG, sparse methods, Event Related Potential, Mixed norm.

1. Introduction

Brain Computer Interfaces (BCI) are systems that help disabled people communicating with their environment through the use of brain signals [14]. At the present time, one of the most prominent BCI is based on electroencephalography (EEG) because of its low-cost, portability and its non-invasiveness. Among EEG based BCI, a paradigm of interest is the one based on event-related potentials (ERP) which are responses of the brain to some external stimuli. In this context, the innermost part of a BCI is the pattern recognition stage which has to correctly recognize presence of these ERPs. However, EEG signals are blurred due to the diffusion of the skull and the skin [27]. Furthermore, EEG recordings are highly contaminated by noise of bi-

*Corresponding author

Email addresses: `remi.flamary@unice.fr` (R. Flamary), `alain.rakoto@insa-rouen.fr` (A. Rakotomamonjy)

ological, instrumental and environmental origins. For addressing these issues, advanced signal processing and machine learning techniques have been employed to learn ERP patterns from training EEG signals leading to robust systems able to recognize the presence of these events [8, 31, 7, 26, 18, 32]. Note that while some ERPs are used for generating BCI commands, some others can be used for improving BCI efficiency. Indeed, recent studies have also tried to develop algorithms for automated recognition of error-related potentials [16]. These potentials are responses elicited when a subject commits an error in a BCI task or observes an error [17, 9] and thus they can help in correcting errors or in providing feedbacks to BCI user's.

In this context of automated recognition of event-related potentials for BCI systems, reducing the number of EEG sensors is of primary importance since it reduces the implementation cost of the BCI by minimizing the number of EEG sensor, and speeding up experimental setup and calibration time. For this purpose, some studies have proposed to choose relevant sensors according to prior knowledge of brain functions. For instance, sensors located above the motor cortex region are preferred for motor imagery tasks and while for visual Event Related Potential (ERP), sensors located on the visual cortex are favored [22]. Recent works have focused on automatic sensor selection adapted to the specificity of a subject [19, 23, 35, 31, 10, 21]. For instance, Rakotomamonjy et al. [30] performed a recursive backward sensor selection using cross-validation classification performances as an elimination criterion. Another approach for exploring subset sensors have been proposed by [35], it consists in using a genetic algorithm for sensor selection coupled with an artificial neural networks for prediction. Those methods has been proven efficient but computationally demanding. A quicker way is to estimate the relevance of the sensors in terms of Signal to Noise Ratio (SNR) [31] and to keep the most relevant ones. Note that this approach does not optimize a discrimination criterion, although the final aim is a classification task. Recently, van Gerven et al. [34] proposed a graceful approach for embedding sensor selection into a discriminative framework. They performed sensor selection and learn a decision function by solving a unique optimization problem. In their framework, a logistic regression classifier is learned and the group-lasso regularization, also known as $\ell_1 - \ell_2$ mixed-norm, is used to promote sensor selection. They have also investigated the use of this groupwise regularization for frequency band selection and their applications to transfer learning. The same idea has been explored by Tomioka et al. [33] which also considered groupwise regularization for classifying EEG signals. In this work, we go beyond these studies by providing an in-depth study of the use of mixed-norms for sensor selection in a single subject setting and by discussing the utility of mixed-norms when learning decision functions for multiple subjects simultaneously.

Our first contribution addresses the problem of robust sensor selection embedded into a discriminative framework. We broaden the analysis of van Gerven et al. [34] by considering regularizers which forms are $\ell_1 - \ell_q$ mixed-norms, with $(1 \leq q \leq 2)$, as well as adaptive mixed-norms, so as to promote sparsity among group of features or sensors. In addition to providing a sparse and accurate sensor selection, mixed-norm regularization has several advantages. First, sensor selection is cast into an elegant discriminative framework, using for

instance a large margin paradigm, which does not require any additional hyper-parameter to be optimized. Secondly, since sensor selection is jointly learned with the classifier by optimizing an “all-in-one” problem, selected sensors are directed to the goal of discriminating relevant EEG patterns. Hence, mixed-norm regularization helps locating sensors which are relevant for an optimal classification performance.

A common drawback of all the aforementioned sensor selection techniques is that selected set of sensors may vary, more or less substantially, from subject to subject. This variability, is due partly to subject specific differences and partly to acquisition noise and limited number of training examples. In such a case, selecting a robust subset of sensors may become a complex problem. Addressing this issue is the point of our second contribution. We propose a Multi-Task Learning (MTL) framework that helps in learning robust classifiers able to cope with the scarcity of learning examples. MTL is one way of achieving inductive transfer between tasks. The goal of inductive transfer is to leverage additional sources of information to improve the performance of learning on the current task. The main hypothesis underlying MTL is that tasks are related in some ways. In most cases, this relatedness is translated into a prior knowledge, *e.g* a regularization term, that a machine learning algorithm can take advantage of. For instance, regularization terms may promote similarity between all the tasks [15], or enforce classifier parameters to lie in a low dimensional linear subspace [2], or to jointly select the relevant features [29]. MTL has been proven efficient for motor imagery in [1] where several classifiers were learned simultaneously from several BCI subject datasets. Our second contribution is thus focused on the problem of performing sensor selection and learning robust classifiers through the use of an MTL mixed-norm regularization framework. We propose a novel regularizer promoting sensor selection and similarity between classifiers. By doing so, our goal is then to yield sensor selection and robust classifiers that are able to overcome the data scarcity problem by sharing information between the different classifiers to be learned.

The paper is organized as follows. The first part of the paper presents the discriminative framework and the different regularization terms we have considered for channel selection and multi-task learning. The second part is devoted to the description of the datasets, the preprocessing steps applied to each of them and the results achieved in terms of performances and sensor selection. In order to promote reproducible research, the code needed for generating the results in this paper is available of the author’s website ¹.

2. Learning framework

In this section, we introduce our mixed-norm regularization framework that can be used to perform sensor selection in a single task or in a transfer learning setting.

¹URL: <http://remi.flamary.com/soft/soft-gsvm.html>

2.1. Channel selection in a single task learning setting

Typically in BCI problems, one wants to learn a classifier that is able to predict the class of some EEG trials, from a set of learning examples. We denoted as $\{\mathbf{x}_i, y_i\}_{i \in \{1 \dots n\}}$ the learning set such that $\mathbf{x}_i \in \mathbb{R}^d$ is a trial and $y_i \in \{-1, 1\}$ is its corresponding class, usually related to the absence or presence of an event-related potential. In most cases, a trial \mathbf{x}_i is extracted from a multidimensional signal and thus is characterized by r features for each of the p sensors, leading to a dimensionality $d = r \times p$. Our aim is to learn, for a single subject, a linear classifier f that will predict the class of a trial $\mathbf{x} \in \mathbb{R}^d$, by looking at the sign of the function $f(\cdot)$ defined as:

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b \quad (1)$$

with $\mathbf{w} \in \mathbb{R}^d$ the normal vector to the separating hyperplane and $b \in \mathbb{R}$ a bias term. Parameters of this function are learned by solving the optimization problem:

$$\min_{\mathbf{w}, b} \sum_i^n L_o(\mathbf{y}_i, \mathbf{x}_i^T \mathbf{w} + b) + \lambda \Omega(\mathbf{w}) \quad (2)$$

where L_o is a loss function that measures the discrepancy between actual and predicted labels, $\Omega(\cdot)$ a regularization term that expresses some prior knowledge about the learning problem and λ a parameter that balances both terms. In this work, we choose L_o to be the squared hinge loss $L_o(y, \hat{y}) = \max(0, 1 - y\hat{y})^2$, thus promoting a large margin classifier.

2.1.1. Regularization terms

We now discuss different regularization terms that may be used for single task learning along with their significances in terms of channel selection.

ℓ_2 norm. The first regularization term that comes to mind is the standard squared ℓ_2 norm regularization:

$$\Omega_2(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2 \quad (3)$$

where $\|\cdot\|_2$ is the Euclidean norm. This is the common regularization term used for SVMs and it will be considered in our experiments as the baseline approach. Intuitively, this regularizer tends to downweight the amplitude of each component of \mathbf{w} leading to a better control of the margin width of our large-margin classifier and thus it helps in reducing overfitting.

ℓ_1 norm. When only few of the features are discriminative for a classification task, a common way to select the relevant ones is to use an ℓ_1 norm of the form

$$\Omega_1(\mathbf{w}) = \sum_{i=1}^d |w_i| \quad (4)$$

as a regularizer [4]. Owing to its mathematical properties (non-differentiability at 0), unlike the ℓ_2 norm, this regularization term promotes sparsity, which means that at optimality of problem (2), some components of \mathbf{w} are exactly 0. In a Bayesian framework, the ℓ_1 norm is related to the use of prior on \mathbf{w} that forces its component to vanish [34]. This is typically obtained by means of Laplacian prior over the weight. However, ℓ_1 norm ignores the structure of the features (which may be grouped by sensors) since each component of \mathbf{w} is treated independently to the others yielding thus to feature selection but not to sensor selection.

$\ell_1 - \ell_q$ mixed-norm. A way to take into account the fact that features are structured, is to use a mixed-norm that will group them and regularize them together. Here, we consider mixed-norm of the form

$$\Omega_{1-q}(\mathbf{w}) = \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_q \quad (5)$$

with $1 \leq q \leq 2$ and \mathcal{G} being a partition of the set $\{1, \dots, d\}$. Intuitively, this $\ell_1 - \ell_q$ mixed-norm can be interpreted as an ℓ_1 norm applied to the vector containing the ℓ_q norm of each group of features. It promotes sparsity on each \mathbf{w}_g norm and consequently on the \mathbf{w}_g components as well. For our BCI problem, a natural choice for \mathcal{G} is to group the features by sensors yielding thus to p groups (one per sensor) of r features as reported in Figure 1. Note that unlike the $\ell_1 - \ell_2$ norm as used by van Gerven et al. [34] and Tomioka et al. [33], the use of an inner ℓ_q norm leads to more flexibility as it spans from the $\ell_1 - \ell_1$ (equivalent to the ℓ_1 -norm and leading thus to unstructured feature selection) to the $\ell_1 - \ell_2$ which strongly ties together the components of a group. Examples of the use of ℓ_q norm and mixed-norm regularizations in other biomedical contexts can be found for instance in [28, 25].

Adaptive $\ell_1 - \ell_q$. The ℓ_1 and $\ell_1 - \ell_q$ norms described above, are well-known to lead to grouped feature selection. However, they are also known, to lead to poor statistical properties (at least when used with a square loss function) [3]. For instance, they are known to have consistency issue in the sense that, even with an arbitrarily large number of training examples, these norms may be unable to select the true subset of features. In practice, this means that when used in Equation (2), the optimal weight vector \mathbf{w} will tend to over-estimate the number of relevant sensors. These issues can be addressed by considering an adaptive $\ell_1 - \ell_q$ mixed-norm of the form [36, 3]:

$$\Omega_{a:1-q}(\mathbf{w}) = \sum_{g \in \mathcal{G}} \beta_g \|\mathbf{w}_g\|_q \quad (6)$$

where the weights β_g are selected so as to enhance the sparsity pattern of \mathbf{w} . In our experiments, we obtain them by first solving the $\ell_1 - \ell_q$ problem with $\beta_g = 1$, which outputs an optimal parameter \mathbf{w}^* , and by finally defining $\beta_g = 1/\|\mathbf{w}_g^*\|_q$. Then, solving the weighted $\ell_1 - \ell_q$ problem yields an optimal solution with increased sparsity pattern compared to \mathbf{w}^* since the β_g augments the penalization of groups with norm $\|\mathbf{w}_g^*\|_q$ smaller than 1.

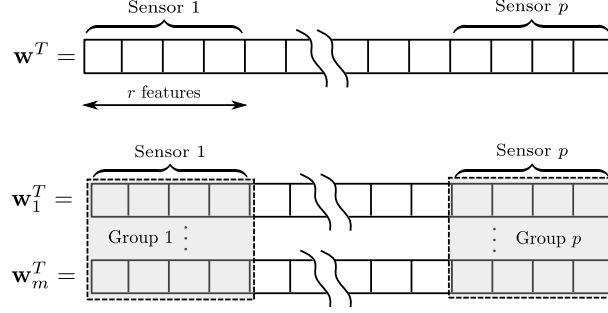


Figure 1: Examples of feature grouping for (top) single task and (bottom) multiple task learning.

2.1.2. Algorithms

Let us now discuss how problem (2) is solved when one of these regularizers is in play.

Using the ℓ_2 norm regularization makes the problem differentiable. Hence a first or second-order descent based algorithm can be considered [11].

Because the other regularizers are not differentiable, we have deployed an algorithm [12] tailored for minimizing objective function of the form $f_1(\mathbf{w}) + f_2(\mathbf{w})$ with f_1 a smooth and differentiable convex function with Lipschitz constant L and f_2 a continuous and convex non-differentiable function having a simple proximal operator, *i.e.* a closed-form or an easy-to-compute solution of the problem:

$$\text{prox}_{f_2}(\mathbf{v}) := \underset{\mathbf{u}}{\text{argmin}} \frac{1}{2} \|\mathbf{v} - \mathbf{u}\|_2^2 + f_2(\mathbf{u}) \quad (7)$$

Such an algorithm, known as forward-backward splitting [12] is simply based on the following iterative approach,

$$\mathbf{w}^{k+1} = \text{prox}_{\frac{1}{\gamma} f_2}(\mathbf{w}^k - \gamma \nabla_{\mathbf{w}} f_1(\mathbf{w}^k)) \quad (8)$$

with γ being a stepsize in the gradient descent. This algorithm can be easily derived by considering, instead of directly minimizing $f_1(\mathbf{w}) + f_2(\mathbf{w})$, an iterative scheme which at each iteration replace f_1 with a quadratic approximation of $f_1(\cdot)$ in the neighborhood of \mathbf{w}^k . Hence, \mathbf{w}^{k+1} is the minimizer of :

$$f_1(\mathbf{w}^k) + \langle \nabla_{\mathbf{w}} f_1(\mathbf{w}^k), \mathbf{w} - \mathbf{w}^k \rangle + \frac{\gamma}{2} \|\mathbf{w} - \mathbf{w}^k\|_2^2 + f_2(\mathbf{w})$$

which closed-form is given in Equation (8). This algorithm is known to converge towards a minimizer of $f_1(\mathbf{w}) + f_2(\mathbf{w})$ under some weak conditions on the stepsize [12], which is satisfied by choosing for instance $\gamma = \frac{1}{L}$. We can note that the algorithm defined in Equation (8) has the same flavor as a projected gradient algorithm which first, takes a gradient step, and then “projects” back the solution owing to the proximal operator. More details can also be found in [5].

For our problem (2), we choose $f_1(\mathbf{w})$ to be the squared hinge loss and $f_2(\mathbf{w})$ the non-smooth regularizer. The square hinge loss is indeed gradient Lipschitz with a constant L being $2 \sum_{i=1} \|\mathbf{x}_i\|_2^2$. Proof of this statement

is available in Appendix 4.1. Proximal operators of the ℓ_1 and the $\ell_1 - \ell_2$ regularization term can be easily shown to be the soft-thresholding and the block-soft thresholding operator [4]. The general $\ell_1 - \ell_q$ norm does not admit a closed-form solution, but its proximal operator can be simply computed by means of an iterative algorithm [29]. More details on these proximal operators are also available in Appendix 4.3.

2.2. Channel selection and transfer learning in multiple task setting

We now address the problem of channel selection in cases where training examples for several subjects are at our disposal. We have claimed that in such a situation, it would be beneficial to learn the decision functions related to all subjects simultaneously, while inducing selected channels to be alike for all subjects, as well as inducing decision function parameters to be related in some sense. These two hypotheses make reasonable sense since brain regions related to the appearance of a given ERP are expected to be somewhat location-invariant across subjects. For solving this problem, we apply a machine learning paradigm, known as multi-task learning, where in our case, each task is related to the decision function of a given subject and where the regularizer should reflect the above-described prior knowledge on the problem. Given m subjects, the resulting optimization problem boils down to be

$$\min_{\mathbf{W}, \mathbf{b}} \sum_t^m \sum_{i=1}^{n_t} L(y_{i,t}, \mathbf{x}_{i,t}^T \mathbf{w}_t + \mathbf{b}_t) + \Omega_{\text{mtl}}(\mathbf{W}) \quad (9)$$

with $\{\mathbf{x}_{i,t}, y_{i,t}\}_{i \in \{1 \dots n_t\}}$ being the training examples related to each task $t \in 1 \dots m$, $(\mathbf{w}_t, \mathbf{b}_t)$ being the classifier parameters for task t and $\mathbf{W} = [\mathbf{w}_1 \dots \mathbf{w}_m] \in \mathbb{R}^{d \times m}$ being a matrix concatenating all vectors $\{\mathbf{w}_t\}$. Note that the multi-task learning framework applied to single EEG trial classification have already been investigated by van Gerven et al. [34]. The main contribution we bring compared to their works is the use of regularizer that explicitly induces all subject classifiers to be similar to an average one, in addition to a regularizer that enforces selected channels to be the same for all subjects. The intuition behind this point is : we believe that since the classification tasks we are dealing with, are similar for all subjects and all related to the same BCI paradigm, selected channels and classifier parameters should not differ that much from subject to subject. We also think that inducing task parameters to be similar may be more important than enforcing selected channels to be similar when the number of training examples is small since it helps in reducing overfitting. For this purpose, we have proposed a novel regularization term of the form :

$$\Omega_{\text{mtl}}(\mathbf{W}) = \lambda_r \sum_{g \in \mathcal{G}'} \|\mathbf{W}_g\|_2 + \lambda_s \sum_{t=1}^m \|\mathbf{w}_t - \hat{\mathbf{w}}\|_2^2 \quad (10)$$

where $\hat{\mathbf{w}} = \frac{1}{m} \sum_t \mathbf{w}_t$ is the average classifier across tasks and \mathcal{G}' contains non-overlapping groups of components from matrix \mathbf{W} . The first term in Equation (10) is a mixed-norm term that promotes group regularization. In this work, we defined groups in \mathcal{G}' based on the sensors, which means that all the features across subject related to a given sensor are in the same group g , leading to p groups of $r \times m$ feature, as depicted

in Figure 1. The second term is a similarity promoting term as introduced in Evgeniou et al. [15]. It can be interpreted as a term enforcing the minimization of the classifier’s parameter variance. In other words, it promotes classifiers to be similar to the average one, and it helps improving performances when the number of learning examples for each task is limited, by reducing over-fitting. Note that λ_r and λ_s respectively control the sparsity induced by the first term and the similarity induced by the second one. Hence, when setting $\lambda_s = 0$, the regularizer given in Equation (10) boils down to be similar to the one used by van Gerven et al. [34]. Note that in practice λ_r and λ_s are selected by means of a nested cross-validation which aims at classification accuracy. Thus, it may occur that classifier similarity is preferred over sensor selection leading to robust classifiers which still use most of the sensors.

Similarly to the single task optimization framework given in Equation (2), the objective function for problem (9) can be expressed as a sum of gradient Lipschitz continuous term $f_1(\mathbf{W}) = \sum_{t,i}^{m,n} L(\cdot) + \lambda_s \sum_{t=1}^m \|\mathbf{w}_t - \hat{\mathbf{w}}\|_2^2$ and a non-differentiable term $f_2(\mathbf{W}) = \lambda_r \sum_{g \in \mathcal{G}'} \|\mathbf{W}_g\|_2$ having a closed-form proximal operator (see Appendix 4.2). Hence, we have again considered a forward-backward splitting algorithm which iterates are given in Equation (8).

3. Numerical experiments

We now present how these novel approaches perform on different BCI problems. Before delving into the details of the results, we introduce the simulated and real datasets.

3.1. Experimental Data

We have first evaluated the proposed approaches on a simple simulated P300 dataset generated as follows. A P300 wave is extracted using the grand average of a single subject data from the EPFL dataset described in the following. We generate 11000 simulated examples with 8 discriminative channels containing the P300 out of 16 channels for positive examples. A Gaussian noise of standard deviation 0.2 is added to all signals making the dataset more realistic. 1000 of these examples have been used for training.

The first real P300 dataset we used is the EPFL dataset, based on eight subjects performing P300 related tasks [20]. The subjects were asked to focus on one of the $3 \times 2 = 6$ images on the screen while the one of the images is flashed at random. The EEG signals were acquired from 32 channels, sampled at 1024 Hz and 4 recording sessions per subject have been realized. Signals are pre-processed exactly according to the steps described in [20] : a $[1, 8]$ Hz bandpass Butterworth filter of order 3 is applied to all signals followed by a downsampling. Hence, for each trial (training example), we have 8 time-sample features per channel corresponding to a 1000 ms time-window after stimulus, which leads to 256 features for all channels ($32 \times 8 = 256$ features). On the overall, the training set of a given subject is composed of about 3000 trials.

Another P300 dataset, recorded by the Neuroimaging Laboratory of Universidad Autónoma Metropolitana (UAM, Mexico) [24], has also been utilized. The data have been obtained from 30 subjects performing P300 spelling tasks on a 6×6 virtual keyboard. Signals are recorded over 10 channels leading thus to a very challenging dataset for sensor selection, as there are just few sensors left to select. For this dataset, we only use the first 3 sessions in order to have the same number of trials for all subjects (≈ 4000 samples). The EEG signals have been pre-processed according to the following steps : a $[2, 20]$ Hz Chebychef bandpass filter of order 5 is first applied followed by a decimation, resulting in a post-stimulus time-window of 31 samples per channels. Hence, each trial is composed of 310 (10×31) features .

We have also studied the effectiveness of our methods on an Error Related Potential (ErrP) dataset that has been recorded in the GIPSA Lab. The subjects were asked to memorize the position of 2 to 9 digits and to remind the position of one of these digits, operation has been repeated 72 times for each subject. The signal following the visualization of the result (correct/error on the memorized position) was recorded from 31 electrodes and sampled at 512 Hz. Similarly to Jrad et al. [21], a $[1, 10]$ Hz Butterworth filter of order 4 and a downsampling has been applied to all channel signals. Finally, a time window of 1000ms is considered as a trial (training example) with a dimensionality of $16 \times 31 = 496$.

3.2. Evaluation criterion, methods and experimental protocol

We have compared several regularizers that induce feature/channel selection embedded in the learning algorithm, in a single subject learning setting as defined in Equation (2). The performance measure commonly used in BCI Competitions [8] is the Area Under the Roc Curve (AUC). This measure is an estimate of the probability for a positive class to have a higher score than a negative class. It makes particularly sense to use AUC when evaluating a P300 speller as the letter in the keyboard is usually chosen by comparing score returned by the classifier for every column or line. In addition, AUC does not depend on the proportion of positive/negative examples in the data which makes it more robust than classification error rate. Our baseline algorithm is an SVM, which uses an ℓ_2 regularizer and thus does not perform any selection. Using an ℓ_1 regularizer yields a classifier which embeds feature selection, denoted as SVM-1 in the sequel. Three mixed-norm regularizers inducing sensor selection have also been considered : an $\ell_1 - \ell_2$ denoted as GSVM-2, and $\ell_1 - \ell_q$ referred as GSVM-q, with q being selected in the set $\{1, 1.2, \dots, 1.8, 2\}$ by a nested cross-validation stage, and adaptive $\ell_1 - \ell_q$ norm, with $q = 2$ denoted as GSVM-a.

For the multi-task learning setting, two MTL methods were compared to two baseline approaches which use all features, namely a method that treats each tasks separately by learning one SVM per task (SVM), and a method denoted as SVM-Full, which on the contrary learns an unique SVM from all subject datasets. The two MTL methods are respectively a MTL as described in Equation (9), denoted as MGSVM-2s and the same MTL but without similarity-promoting regularization term, which actually means that we set $\lambda_s = 0$,

Methods	Avg AUC	AUC p-val	Avg Sel	F-measure
SVM	79.79	-	100.00	66.67
GSVM-1	79.32	0.027	98.75	67.25
GSVM-2	80.96	0.004	62.50	89.72
GSVM-p	80.74	0.020	63.12	89.40
GSVM-a	80.51	0.014	45.62	93.98

Table 1: Performance results on the simulated dataset : the average performance in AUC (in %), the average percent of selected sensors (Sel) and the F-measure of the selected channels (in %). Best results for each performance measure are in bold. The p-value refers to the one of a Wilcoxon signrank test with SVM as a baseline.

indicated as MGSVM-2. For these approaches, performances are evaluated as the average AUC of the decision functions over all the subjects.

The experimental setup is described in the following. For each subject, the dataset is randomly split into a training set of $n = 1000$ trials and a test set containing the rest of the trials. The regularization parameter λ has been selected from a log-spaced grid ($[10^{-3}, 10^1]$) according to a nested 3-fold cross-validation step on the training set. When necessary, the selection of q is also included in this CV procedure. Finally, the selected value of λ is used to learn a classifier on the training examples and performances are evaluated on the independent test set. We run this procedure 10 times for every subject and report average performances. A Wilcoxon signed-rank test, which takes ties into account is used to evaluate the statistical difference of the mean performances of all methods compared to the baseline SVM. We believe that such a test is more appropriate for comparing methods than merely looking at the standard deviation due to the high inter-subject variability in BCI problems.

3.3. Results and discussions

We now present the results we achieved on the above-described datasets.

3.3.1. Simulated dataset

Average (over 10 runs) performance of the different regularizers on the simulated dataset are reported in Table 1 through AUC, sensor selection rate and F-measure. This latter criterion measures the relevance of the selected channels compared to the true relevant ones. F-measure is formally defined as

$$\text{F-measure} = 2 \frac{|\mathcal{C} \cap \mathcal{C}^*|}{|\mathcal{C}^*| + |\mathcal{C}|}$$

where \mathcal{C} and \mathcal{C}^* are respectively the set of selected channels and true relevant channels and $|\cdot|$ here denotes the cardinality of a set. Note that if the selected channels are all the relevant ones, then the F-measure is equal to one. Most of the approaches provide similar AUC performances. We can although highlight that group-regularization approaches (GSVM-2, GSVM-p, GSVM-a) drastically reduce the number of selected channels

Datasets	EPFL Dataset (8 Sub., 32 Ch.)			UAM Dataset (30 Sub., 10 Ch.)			ErrP Dataset (8 Sub., 32 Ch)		
Methods	Avg AUC	Avg Sel	p-value	Avg AUC	Avg Sel	p-value	Avg AUC	Avg Sel	p-value
SVM	80.35	100.00	-	84.47	100.00	-	76.96	100.00	-
SVM-1	79.88	87.66	0.15	84.45	96.27	0.5577	68.84	45.85	0.3125
GSVM-2	80.53	78.24	0.31	84.94	88.77	0.0001	77.29	29.84	0.5469
GSVM-p	80.38	77.81	0.74	84.94	90.80	0.0001	76.84	37.18	0.7422
GSVM-a	79.01	26.60	0.01	84.12	45.07	0.1109	67.25	7.14	0.1484

Table 2: Performance results for the 3 datasets the average performance (over subjects) in AUC (in %), the average percent of selected sensors (Sel) and the p-value of the Wilcoxon signrank test for the AUC when compared to the baseline SVM’s one. Best performing algorithms for each performance measure are in bold.

since only 62% and 45% of the sensors are selected. A clear advantage goes to the adaptive regularization that is both sparser and is more capable of retrieving the true relevant channels.

3.3.2. P300 Datasets

Results for these datasets are reported in Table 2. For the EPFL dataset, all methods achieve performances that are not statistically different. However, we note that GSVM-2 leads to sensor selection (80% of sensor selected) while GSVM-a yields to classifiers that, on average, use 26% of the sensors at the cost of a slight loss in performances (1.5% AUC).

Results for the UAM dataset follow the same trend in term of sensor selection but we also observe that the mixed-norm regularizers yield to increased performances. GSVM-2 performs statistically better than SVM although most of the sensors (9 out of 10) have been kept in the model. This shows that even if few channels have been removed, the group-regularization improves performances by bringing sensor prior knowledge to the problem. We also notice that GSVM-a performance is statistically equivalent to the baseline SVM one while using only half of the sensors and GSVM-p consistently gives similar results to GSVM-2.

To summarize, concerning the performances of the different mixed-norm regularization, we outline that on one hand, GSVM-2 is at worst, equivalent to the baseline SVM while achieving sensor selection and on the other hand GSVM-a yields to a substantial channel selection at the expense of a slight loss of performances.

A visualization of the electrodes selected by GSVM-a can be seen in Figure 2 for the EPFL dataset and in Figure 3 for the UAM dataset. Interestingly, we observe that for the EPFL dataset, the selected channels are highly dependent on the subject. The most recurring ones are : FC1 C3 T7 CP5 P3 PO3 PO4 Pz and the electrodes located above visual cortex O1,Oz and O2. We see sensors from the occipital area that are known to be relevant [22] for P300 recognition, but sensors such as T7 and C3, from other brain regions are also frequently selected. These results are however consistent with those presented in the recent literature [31, 30].

The UAM dataset uses only 10 electrodes that are already known to perform well in P300 recognition problem, but we can see from Figure 3 that the adaptive mixed-norm regularizer further selects some sensors that are

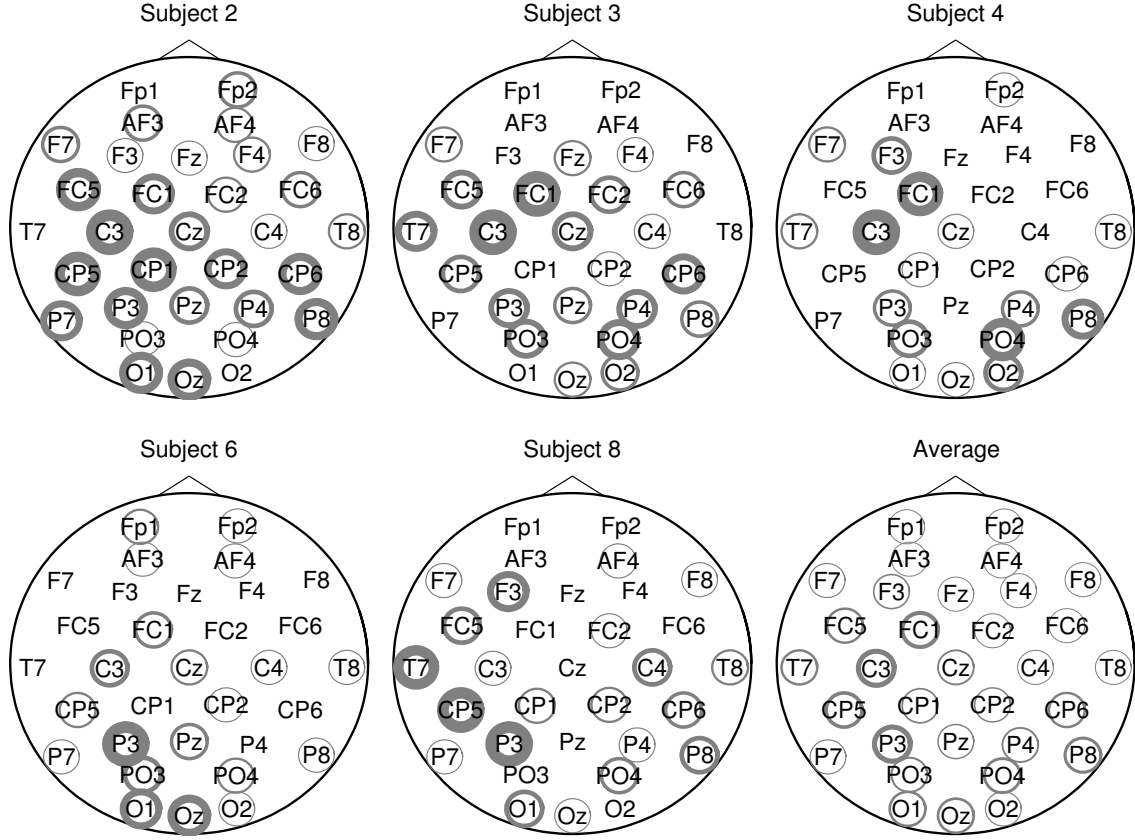


Figure 2: Selected sensors for the EPFL dataset. The line width of the circle is proportional to the number of times the sensor is selected for different splits. No circle means that the sensor has never been selected.

essentially located in the occipital region. Note that despite the good average performances reported in Table 2, some subjects in this dataset achieve very poor performances, of about 50 % of AUC, regardless of the considered method. Selected channels for one of these subjects (Subject 25) are depicted in Figure 3 and interestingly, they strongly differ from those of other subjects providing rationales for the poor AUC.

We have also investigated the impact of sparsity on the overall performance of the classifiers. To this aim, we have plotted the average performance of the different classifiers as a function of the number of selected sensors. These plots are depicted in Figure 4 for the EPFL dataset and on Figure 5 for the UAM dataset. For both datasets, GSVM-a frequently achieves a better AUC for a given level of sparsity. For most of the subjects, GSVM-a performs as well as SVM but using far less sensors. A rationale may be that, in addition to selecting the relevant sensors, GSVM-a may provide a better estimation of the classifier parameters leading to better performances for a fixed number of sensors. As a summary, we suggest thus the use of an adaptive mixed-norm regularizer instead of an $\ell_1 - \ell_2$ mixed-norm as in van Gerven et al. [34] when sparsity and channel selection is of primary importance.

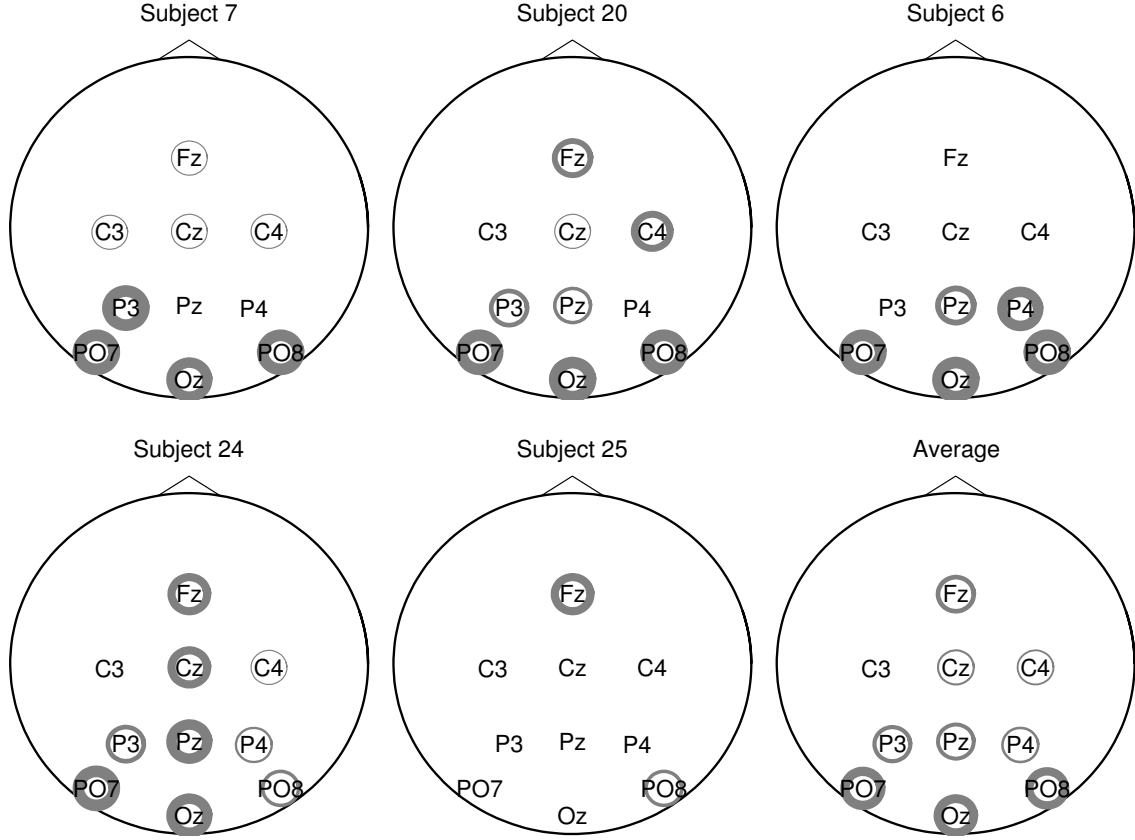


Figure 3: Selected sensors for the UAM dataset. The line width of the circle is proportional to the number of times the sensor is selected for different splits. No circle means that the sensor has never been selected.

3.3.3. ErrP Dataset

The ErrP dataset differs from the others as its number of examples is small (72 examples per subject). The same experimental protocol as above has been used for evaluating the methods but only 57 examples out of 72 have been retained for validation/training. Classification performances are reported on Table 2. For this dataset, the best performance is achieved by GSVM-2 but the Wilcoxon test shows that all methods are actually statistically equivalent. Interestingly, many channels of this dataset seem to be irrelevant for the classification task. Indeed, GSVM-2 selects only 30% of them while GSVM-a uses only 7% of the channels at the cost of 10% AUC loss. We believe that this loss is essentially caused by the aggressive regularization of GSVM-a and the difficulty to select the regularization parameter λ using only a subset of the 57 training examples. Channels selected by GSVM-2 can be visualized on Figure 6. Despite the high variance in terms of selected sensors, probably due to the small number of examples, sensors in the central area seem to be the most selected one, which is consistent with previous results in ErrP [13].

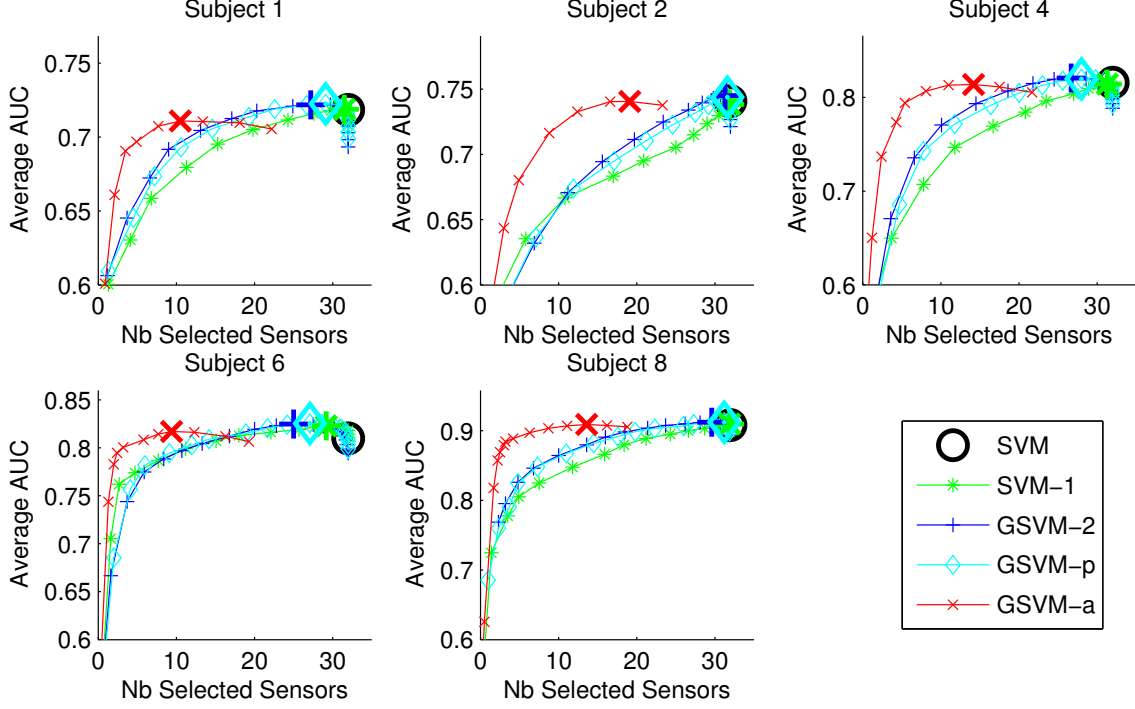


Figure 4: Performance vs sensor selection visualisation for the EPFL dataset. The large marker corresponds to the best model along the regularization path.

3.3.4. Multi-task Learning

We now evaluate the impact of the approach we proposed in Equation (9) and (10) on the P300 datasets. We expect that since multi-task learning allows to transfer some information between the different classification tasks, it will help in leveraging classification performances especially when the number of available training examples is small. Note that the ErrP dataset has not been tested in this MTL framework, because the above-described results suggest an important variance in the selected channels for all subjects. Hence, we believe that this learning problem does not fit into the prior knowledge considered through Equation (10).

We have followed the same experimental protocol as for the single task learning except that training and test sets have been formed as follows. We first create training and test examples for a given subject by randomly splitting all examples of that subject, and then gather all subject’s training/test sets to form the multi-task learning training/test sets. Hence, all the subjects are equally represented in these sets. A 3-fold nested cross-validation method is performed in order to automatically select the regularization terms (λ_r and λ_s).

Performances of the different methods have been evaluated for increasing number of training examples per subject and are reported in Figure 7. We can first see that for the EPFL dataset, MGSVM-2 and MGSVM-2s yield a slight but consistent improvement over the single-task classifiers (SVM-Full being a single classifier trained on all subject’s examples and SVM being the average performances of subject-specific classifiers).

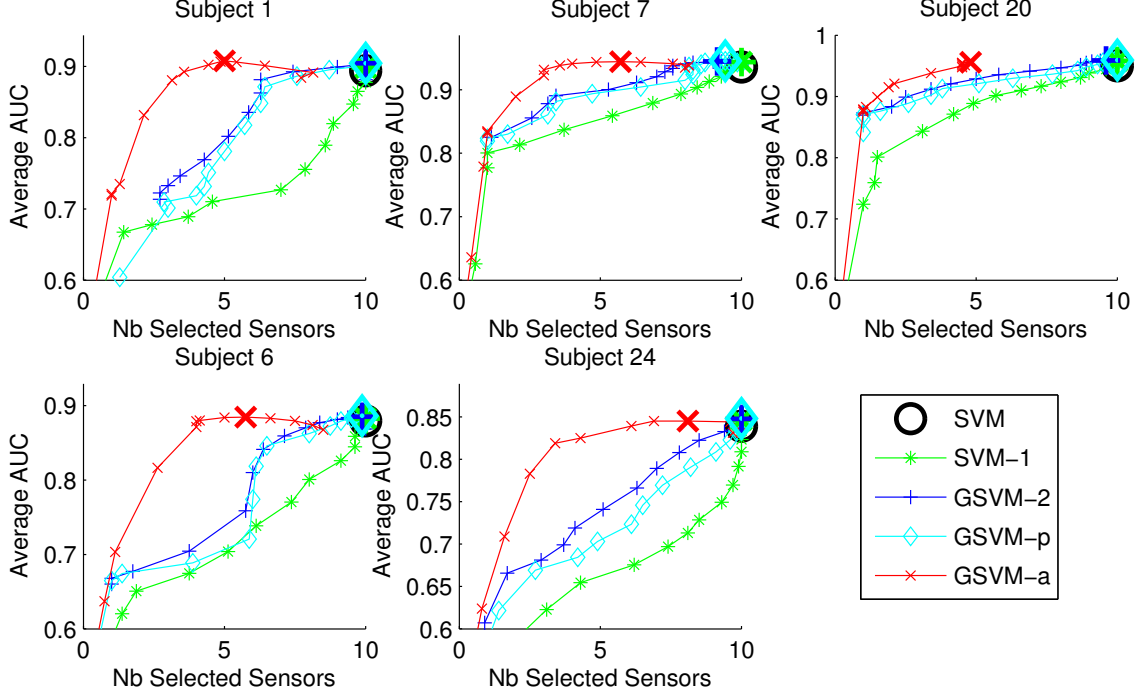


Figure 5: Performance vs sensor selection visualisation for the UAM dataset. The large marker corresponds to the best model along the regularization path.

The poor performances of the SVM-Full approach is probably due to the high inter-subject variability in this dataset, which includes impaired patients.

For the UAM dataset, results are quite different since the SVM-Full and MGSVM-2s shows a significant improvement over the single-task learning. We also note that, when only the joint channel selection regularizer is in play (MGSVM-2), multi-task learning leads to poorer performance than the SVM-Full for a number of trials lower than 500. We justify this by the difficulty of achieving appropriate channel selection based only on few training examples, as confirmed by the performance of GSVM-2. From Figure 8, we can see that the good performance of MGSVM-2s is the outcome of performance improvement of about 10% AUC over SVM, achieved on some subjects that perform poorly. More importantly, while performances of these subjects are significantly increased, those that performs well still achieve good AUC scores. In addition, we emphasize that these improvements are essentially due to the similarity-inducing regularizer.

For both datasets, the MTL approach MGSVM-2s is consistently better than those of other single-task approaches thanks to the regularization parameters λ_r and λ_s that can adapt to the inter-subject similarity (weak similarity for EPFL and strong similarity for UAM). These are interesting results showing that multi-task learning can be a way to handle the problem related to some subjects that achieve poor performances. Moreover, results also indicate that multi-task learning is useful for drastically shortening the calibration time. For instance, for the UAM dataset, 80% AUC was achieved using only 100 training examples (less than 1

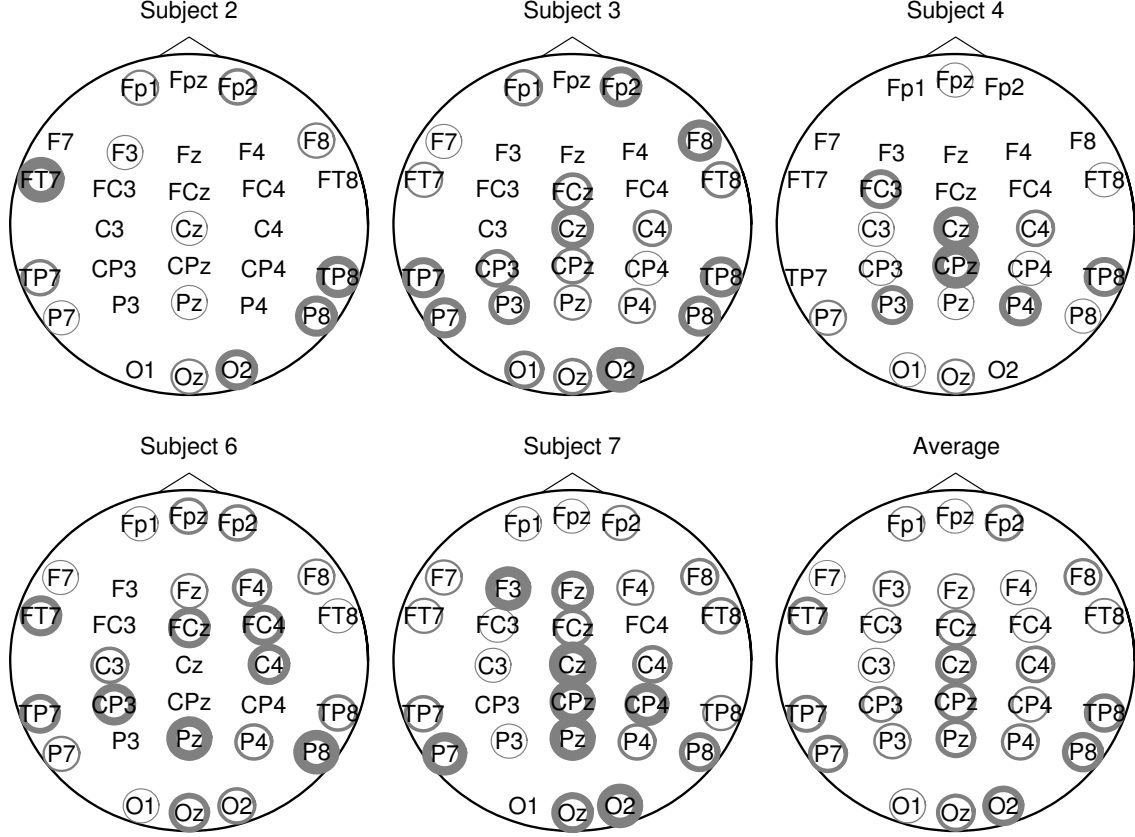


Figure 6: Selected Sensors for the ERP dataset. The line width of the circle is proportional to the number of times the sensor is selected. No circle means that the sensor has never been selected.

minute of training example recordings). Note that the validation procedure tends to maximize performances, and does not lead to sparse classifiers for MTL approaches. As shown in Figures 2 and 3, the relevant sensors are quite different between subjects thus a joint sensor selection can lead to a slight loss of performances, hence the tendency of the cross-validation procedure to select non-sparse classifiers.

4. Conclusion

In this work, we have investigated the use of mixed-norm regularizers for discriminating Event-Related Potentials in BCI. We have extended the discriminative framework of van Gerven et al. [34] by studying general mixed-norms and proposed the use of the adaptive mixed-norms as sparsity-inducing regularizers. This discriminative framework has been broadened to the multi-task learning framework where classifiers related to different subjects are jointly trained. For this framework, we have introduced a novel regularizer that induces channel selection and classifier similarities. The different proposed approaches were tested on three different datasets involving a substantial number of subjects. Results from these experiments have highlighted that the $\ell_1 - \ell_2$ regularizer has been proven interesting for improving classification performance whereas adaptive

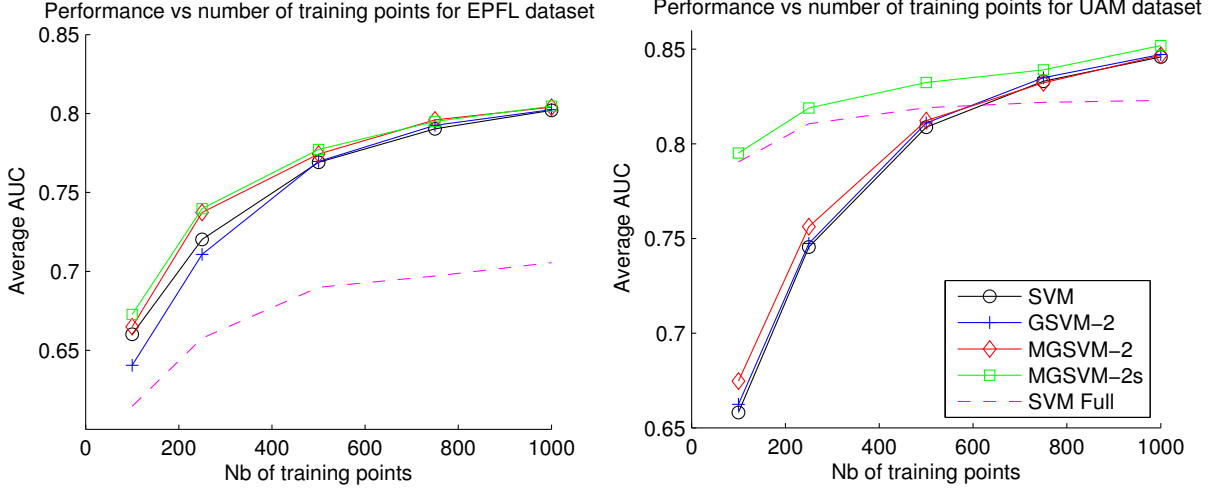


Figure 7: Multi-task learning performances (AUC) for the EPFL (left plot) and UAM (right plot) datasets for different number of training examples per subject.

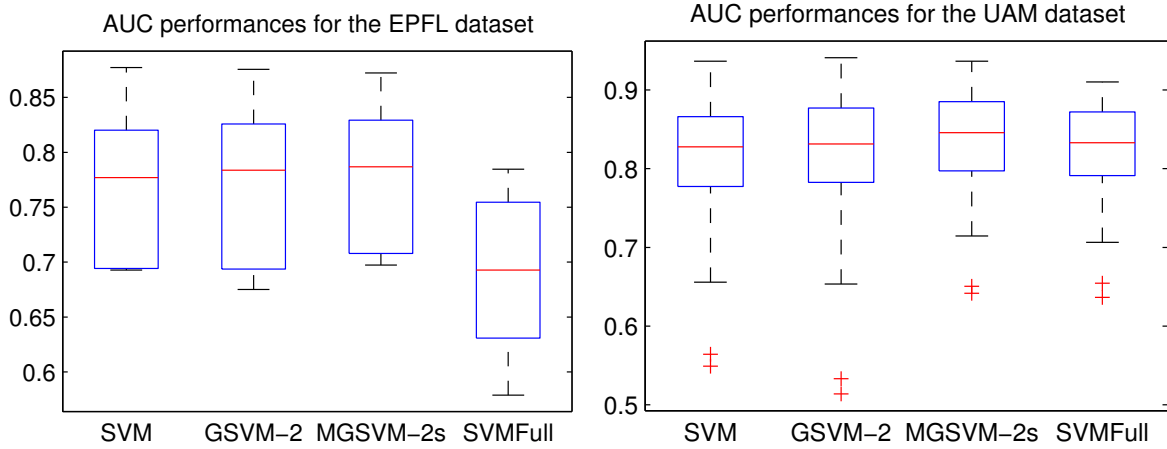


Figure 8: AUC performances comparison with EPFL (left plot) and UAM (right plot) for 500 training examples per subject.

mixed-norm is the regularizer to be considered when sensor selection is the primary objective. Regarding the multi-task learning framework, our most interesting finding is that this learning framework allows, by learning more robust classifiers, significant performance improvement on some subjects that perform poorly in a single-task learning context.

In future work, we plan to investigate a different grouping of the features, such as temporal groups. This kind of group regularization could be for instance used in conjunction with the sensors group in order to promote both feature selection and temporal selection in the classifier. While the resulting problem is still convex, its resolution poses some issues so that a dedicated solver would be necessary.

Another research direction would be to investigate the use of asymmetrical MTL. This could prove handy when a poorly-performing subject will negatively influence the other subject performances in MTL while

improving his own performances. In this case one would like that subject classifier to be similar to the other's classifier without impacting their classifiers.

- [1] Alamgir, M., Grosse-Wentrup, M., Altun, Y., 2010. Multi-task Learning for Brain-Computer Interfaces. In: AI & Statistics.
- [2] Argyriou, A., Evgeniou, T., Pontil, M., 2008. Convex multi-task feature learning. *Machine Learning* 73 (3), 243–272.
- [3] Bach, F., 2008. Consistency of the group Lasso and multiple kernel learning. *Journal of Machine Learning Research* 9, 1179–1225.
- [4] Bach, F., Jenatton, R., Mairal, J., Obozinski, G., 2011. Convex optimization with sparsity-inducing norms. In: Sra, S., Nowozin, S., Wright, S. J. (Eds.), *Optimization for Machine Learning*. MIT Press, pp. 19–53.
- [5] Beck, A., Teboulle, M., 2009. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences* 2, 183–202.
- [6] Bertsekas, D., 1999. *Nonlinear programming*. Athena scientific.
- [7] Blankertz, B., Lemm, S., Treder, M., Haufe, S., Müller, K., 2010. Single-trial analysis and classification of ERP components – a tutorial. *NeuroImage* 56 (2), 814–825.
URL <http://dx.doi.org/10.1016/j.neuroimage.2010.06.048>
- [8] Blankertz, B., Mueller, K.-R., Krusienski, D., Schalk, G., Wolpaw, J., Schloegl, A., Pfurtscheller, G., del R. Millan, J., Schroeder, M., Birbaumer, N., 2006. The BCI competition III: Validating alternative approaches to actual BCI problems. *IEEE Transactions on Neural Systems and Rehabilitation Engineering* 14 (2), 153–159.
- [9] Buttfield, A., Ferrez, P., Millan, J., 2006. Towards a robust bci: Error potentials and online learning. *Neural Systems and Rehabilitation Engineering, IEEE Transactions on* 14 (2), 164–168.
- [10] Cecotti, H., Rivet, B., Congedo, M., Jutten, C., Bertrand, O., Maby, E., Mattout, J., 2011. A robust sensor-selection method for P300 brain-computer interfaces. *Journal of Neural Engineering*.
- [11] Chapelle, O., 2007. Training a support vector machine in the primal. *Neural Comput.* 19 (5), 1155–1178.
- [12] Combettes, P., Pesquet, J., 2011. Proximal splitting methods in signal processing. *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, 185–212.

- [13] Dehaene, S., Posner, M., Tucker, D., 1994. Localization of a neural system for error detection and compensation. *Psychological Science* 5 (5), 303–305.
- [14] Dornhege, G., Millán, J., Hinterberger, T., McFarland, D., Müller, K., 2007. Toward brain-computer interfacing. Vol. 74. Mit Press Cambridge, MA.
- [15] Evgeniou, T., Pontil, M., 2004. Regularized multi-task learning. In: *Proceedings of the tenth Conference on Knowledge Discovery and Data mining*.
- [16] Falkenstein, M., Hohnsbein, J., Hoormann, J., Blanke, L., 1991. Effects of crossmodal divided attention on late erp components. ii. error processing in choice reaction tasks. *Electroencephalography and clinical Neurophysiology* 78 (6), 447–455.
- [17] Ferrez, P., Millán, J., 2007. Error-related eeg potentials in brain-computer interfaces. *Towards Brain-Computer Interfacing*. MIT Press, Cambridge, Massachusetts.
- [18] Gouy-Pailler, C., Congedo, M., Brunner, C., Jutten, C., Pfurtscheller, G., 2010. Nonstationary brain source separation for multiclass motor imagery. *IEEE Trans. on Biomedical Engineering* 57 (2), 469–478.
- [19] Hoffman, U., Yazdani, A., Vesin, J., Ebrahimi, T., 2008. Bayesian feature selection applied in a P300 brain-computer interface. In: *16th European signal processing conference*.
- [20] Hoffmann, U., Vesin, J., Ebrahimi, T., Diserens, K., 2008. An efficient p300-based brain-computer interface for disabled subjects. *Journal of Neuroscience Methods* 167 (1), 115–125.
- [21] Jrad, N., Congedo, M., Phlypo, R., Rousseau, S., Flamary, R., Yger, F., Rakotomamonjy, A., 2011. sw-svm: sensor weighting support vector machines for eeg-based brain-computer interfaces. *Journal of Neural Engineering* 8, 056004.
- [22] Krusienski, D., Sellers, E., McFarland, D., Vaughan, T., Wolpaw, J., 2008. Towards enhanced P300 speller performances. *Journal of neuroscience methods* 167 (1), 15–21.
- [23] Lal, T., Schroder, M., Hinterberger, T., Weston, J., Bogdan, M., Birbaumer, N., Scholkopf, B., 2004. Support vector channel selection in BCI. *Biomedical Engineering, IEEE Transactions on* 51 (6), 1003–1010.
- [24] Ledesma-Ramirez, C., Bojorges Valdez, E., Yáñez Suarez, O., Saavedra, C., Bougrain, L., Gentiletti, G. G., 2010. An Open-Access P300 Speller Database. In: *Fourth International Brain-Computer Interface Meeting*.
- [25] Liu, A., Hao, T., Gao, Z., Su, Y., Yang, Z., 2013. Non-negative mixed-norm convex optimization for mitotic cell detection in phase contrast microscopy. *Computational and Mathematical Methods in Medicine* 2013, 1–10.

- [26] Müller-Gerking, J., Pfurtscheller, G., Flyvbjerg, H., 1999. Designing optimal spatial filters for single-trial EEG classification in a movement task. *Clin. Neurophys.* 110, 787–798.
- [27] Nunez, P. L., Srinivasan, R., 2006. *Electric Fields of the Brain*, 2nd Edition. New York: Oxford Univ Press.
- [28] Rahimi, A., Xu, J., Wang, L., 2013. ℓ_p norm regularization in volumetric imaging of cardiac current sources. *Computational and Mathematical Methods in Medicine* 2013, 1–10.
- [29] Rakotomamonjy, A., Flamary, R., Gasso, G., Canu, S., 2011. lp-lq penalty for sparse linear and sparse multiple kernel multi-task learning,. *IEEE Transactions on Neural Networks* 22 (8), 1307–1320.
- [30] Rakotomamonjy, A., Guigue, V., 2008. BCI competition III: Dataset II - ensemble of SVMs for BCI P300 speller. *IEEE Trans. Biomedical Engineering* 55 (3), 1147–1154.
- [31] Rivet, B., Cecotti, H., Phlypo, R., Bertrand, O., Maby, E., Mattout, J., 2010. EEG sensor selection by sparse spatial filtering in P300 speller brain-computer interface. In: *Engineering in Medicine and Biology Society (EMBC), 2010 Annual International Conference of the IEEE. IEEE*, pp. 5379–5382.
- [32] Salimi-Khorshidi, G., Nasrabadi, A., Golpayegani, M., 2008. Fusion of classic p300 detection methods’ inferences in a framework of fuzzy labels. *Artificial Intelligence in Medicine* 44 (3), 247–259.
- [33] Tomioka, R., Müller, K., 2010. A regularized discriminative framework for EEG analysis with application to brain-computer interface. *NeuroImage* 49 (1), 415–432.
- [34] van Gerven, M., Hesse, C., Jensen, O., Heskes, T., 2009. Interpreting single trial data using groupwise regularisation. *NeuroImage* 46 (3), 665–676.
- [35] Yang, J., Singh, H., Hines, E. L., Schlaghecken, F., Iliescu, D. D., Leeson, M. S., Stocks, N. G., 2012. Channel selection and classification of electroencephalogram signals: An artificial neural network and genetic algorithm-based approach. *Artificial Intelligence in Medicine* 55 (2), 117 – 126.
URL <http://www.sciencedirect.com/science/article/pii/S0933365712000292>
- [36] Zou, H., 2006. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* 101 (476), 1418–1429.

Appendix

4.1. Proof of Lipschitz gradient of the squared Hinge loss

Given the training examples $\{\mathbf{x}_i, y_i\}$, the squared Hinge loss is written as :

$$J = \sum_{i=1}^n \max(0, 1 - y_i \mathbf{x}_i^\top \mathbf{w})^2$$

and its gradient is :

$$\nabla_{\mathbf{w}} J = -2 \sum_i \mathbf{x}_i y_i \max(0, 1 - y_i \mathbf{x}_i^\top \mathbf{w})$$

The squared Hinge loss is gradient Lipschitz if there exists a constant L such that:

$$\|\nabla J(\mathbf{w}_1) - \nabla J(\mathbf{w}_2)\|_2 \leq L \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \quad \forall \mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d.$$

The proof essentially relies on showing that $\mathbf{x}_i y_i \max(0, 1 - y_i \mathbf{x}_i^\top \mathbf{w})$ is Lipschitz itself *i.e* there exists $L' \in \mathbb{R}$ such that

$$\begin{aligned} & \|\mathbf{x}_i y_i \max(0, 1 - y_i \mathbf{x}_i^\top \mathbf{w}_1) - \mathbf{x}_i y_i \max(0, 1 - y_i \mathbf{x}_i^\top \mathbf{w}_2)\| \\ & \leq L' \|\mathbf{w}_1 - \mathbf{w}_2\| \end{aligned}$$

Now let us consider different situations. For a given \mathbf{w}_1 and \mathbf{w}_2 , if $1 - \mathbf{x}_i^\top \mathbf{w}_1 \leq 0$ and $1 - \mathbf{x}_i^\top \mathbf{w}_2 \leq 0$, then the left hand side is equal to 0 and any L' would satisfy the inequality. If $1 - \mathbf{x}_i^\top \mathbf{w}_1 \leq 0$ and $1 - \mathbf{x}_i^\top \mathbf{w}_2 \geq 0$, then the left hand side (lhs) is

$$\begin{aligned} lhs &= \|\mathbf{x}_i\|_2 (1 - \mathbf{x}_i^\top \mathbf{w}_2) \\ &\leq \|\mathbf{x}_i\|_2 (\mathbf{x}_i^\top \mathbf{w}_1 - \mathbf{x}_i^\top \mathbf{w}_2) \\ &\leq \|\mathbf{x}_i\|_2^2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2 \end{aligned} \tag{11}$$

A similar reasoning yields to the same bound when $1 - \mathbf{x}_i^\top \mathbf{w}_1 \geq 0$ $1 - \mathbf{x}_i^\top \mathbf{w}_1 \leq 0$ and $1 - \mathbf{x}_i^\top \mathbf{w}_2 \geq 0$ and $1 - \mathbf{x}_i^\top \mathbf{w}_2 \leq 0$. Thus, $\mathbf{x}_i y_i \max(0, 1 - y_i \mathbf{x}_i^\top \mathbf{w})$ is Lipschitz with a constant $\|\mathbf{x}_i\|^2$. Now, we can conclude the proof by stating that $\nabla_{\mathbf{w}} J$ is Lipschitz as it is a sum of Lipschitz function and the related constant is $\sum_{i=1}^n \|\mathbf{x}_i\|_2^2$.

4.2. Lipschitz gradient for the multi-task learning problem

For the multi-task learning problem, we want to prove that the function

$$\sum_{t=1}^m \sum_{i=1}^n L(y_{i,t}, \mathbf{x}_{i,t}^\top \mathbf{w}_t + \mathbf{b}_t) + \lambda_s \sum_{t=1}^m \|\mathbf{w}_t - \frac{1}{m} \sum_{j=1}^m \mathbf{w}_j\|_2^2$$

is gradient Lipschitz, $L(\cdot, \cdot)$ being the square Hinge loss. From the above results, it is easy to show that the first term is gradient Lipschitz as the sum of gradient Lipschitz functions.

Now, we also show that the similarity term

$$\sum_t \left\| \mathbf{w}_t - \frac{1}{m} \sum_{j=1}^m \mathbf{w}_j \right\|_2^2$$

is also gradient Lipschitz.

This term can be expressed as

$$\begin{aligned} \left\| \mathbf{w}_t - \frac{1}{m} \sum_{j=1}^m \mathbf{w}_j \right\|_2^2 &= \sum_t \langle \mathbf{w}_t, \mathbf{w}_t \rangle - \frac{1}{m} \sum_{i,j=1}^m \langle \mathbf{w}_i, \mathbf{w}_j \rangle \\ &= \mathbf{w}^\top \mathbf{M} \mathbf{w} \end{aligned}$$

where $\mathbf{w}^\top = [\mathbf{w}_1^\top, \dots, \mathbf{w}_m^\top]$ is the vector of all classifier parameters and $\mathbf{M} \in \mathbb{R}^{md \times md}$ is the Hessian matrix of the similarity regularizer of the form

$$\mathbf{M} = \mathbf{I} - \frac{1}{m} \sum_{t=1}^m \mathbf{D}_t$$

with \mathbf{I} the identity matrix and \mathbf{D}_t a block matrix with \mathbf{D}_t a $(t-1)$ -diagonal matrix where each block is an identity matrix \mathbf{I} with appropriate circular shift. \mathbf{D}_t is thus a $(t-1)$ row-shifted version of \mathbf{I} .

Once we have this formulation, we can use the fact that a function f is gradient Lipschitz of constant L if the largest eigenvalue of its Hessian is bounded by L on its domain [6]. Hence, since we have

$$\|\mathbf{M}\|_2 \leq \|\mathbf{I}\|_2 + \frac{1}{m} \sum_{t=1}^m \|\mathbf{D}_t\|_2 = 2$$

the Hessian matrix of the similarity term $2 \cdot \mathbf{M}$ has consequently bounded eigenvalues. This concludes the proof that the function $\mathbf{w}^\top \mathbf{M} \mathbf{w}$ is gradient Lipschitz continuous.

4.3. Proximal operators

4.3.1. ℓ_1 norm

the proximal operator of the ℓ_1 norm is defined as :

$$\text{prox}_{\lambda \|\mathbf{x}\|_1}(\mathbf{u}) = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

and has the following closed-form solution for which each component is

$$[\text{prox}_{\lambda \|\mathbf{x}\|_1}(\mathbf{u})]_i = \text{sign}(u_i)(|u_i| - \lambda)_+$$

4.3.2. $\ell_1 - \ell_2$ norm

the proximal operator of the $\ell_1 - \ell_2$ norm is defined as :

$$\text{prox}_{\lambda \sum_{g \in \mathcal{G}} \|\mathbf{x}_g\|_2}(\mathbf{u}) = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_2^2 + \lambda \sum_{g \in \mathcal{G}} \|\mathbf{x}_g\|_2$$

the minimization problem can be decomposed into several ones since the indices g are separable. Hence, we can just focus on the problem

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_2^2 + \lambda \|\mathbf{x}\|_2$$

which minimizer is

$$\begin{cases} 0 & \text{if } \|\mathbf{u}\|_2 \leq \lambda \\ (1 - \frac{\lambda}{\|\mathbf{u}\|_2})\mathbf{u} & \text{otherwise} \end{cases}$$